

A Survey: Deepfake and Current Technologies for Solutions

Sayan Banerjee¹, Sumit Kumar Yadav¹, Ankit Dhara¹ and Md Ajjij^{1,*},[†]

¹Department of Computer Science and Technology, University of North Bengal, Raja Rammohunpur, Darjeeling, West Bengal, 734013, India

Abstract

This paper offers a detailed survey of deepfake detection methods, addressing the challenges posed by the fast-paced advancements in deepfake technology. It provides an overview of various detection techniques, examining their effectiveness in identifying manipulated content. The survey covers traditional detection strategies, such as digital forensics and watermarking, as well as modern AI-driven approaches like convolutional and recurrent neural networks. The study delves into the key features of deepfake technology, which leverages advanced machine learning models, particularly Generative Adversarial Networks (GANs), to manipulate video, audio, and images. These techniques have led to the creation of highly realistic synthetic media that is increasingly difficult to detect, raising serious concerns about privacy, misinformation, and security. Recent progress in deepfake detection has focused on improving the accuracy and efficiency of real-time solutions. Approaches that integrate visual, audio, and behavioural cues have demonstrated significant potential in distinguishing authentic content from fake media. Despite these advancements, there remains an urgent need for detection systems that can generalize effectively across different types of deepfakes, as many current models struggle with previously unseen or extremely realistic synthetic content. The survey reviews a broad spectrum of detection methods, assessing their strengths, weaknesses, and performance on various datasets. It also identifies gaps in the current research landscape and suggests directions for future work, emphasizing the importance of developing more robust and scalable detection frameworks.

Keywords

Deepfake, Survey, Advanced machine learning models, Generative Adversarial Networks (GANs), Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN)

1. Introduction

Deepfakes, a term combining "deep learning" and "fake", describe highly convincing synthetic media produced using advanced machine learning techniques. Emerging in 2017, deepfakes initially focused on facial manipulation in videos. Since then, the technology has expanded to encompass audio and image alteration. Using algorithms like Generative Adversarial Networks (GANs), deepfakes can realistically swap faces, modify facial expressions, and even mimic voices, making it increasingly challenging to distinguish between genuine and synthetic content. Although initially developed for entertainment purposes, deepfake technology has evolved rapidly, bringing with it significant implications for digital privacy, security, and the reliability of online information.

The swift advancement of deepfake technology is both impressive and concerning. As the algorithms become more sophisticated, so does the quality of synthetic content. This progress has sparked worries about the potential misuse of deepfakes for spreading misinformation, committing fraud, and facilitating identity theft. Deepfakes have already been used in disinformation campaigns, influencing public perception and casting doubt on media authenticity. Their potential to erode trust in individuals and institutions underscores the urgent need for effective detection and prevention measures.

This paper seeks to offer an in-depth survey of the existing methods for detecting and mitigating deepfakes. By examining various techniques, such as facial feature analysis, biometric inconsistencies,

The 2024 Sixth Doctoral Symposium on Intelligence Enabled Research (DoSIER 2024), November 28–29, 2024, Jalpaiguri, India

*Corresponding author.

[†]These authors contributed equally.

✉ banerjeesayan554@gmail.com (S. Banerjee); sk9373279@gmail.com (S. K. Yadav); ankitdhara8250@gmail.com (A. Dhara); mdajij@nbu.ac.in (M. Ajjij)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

and behavioural patterns, the study assesses the effectiveness of these approaches across different datasets and scenarios. The goal is to highlight current solutions while identifying research gaps and suggesting future directions to address the growing sophistication of deepfake technology.

The motivation for this survey stems from the increasing need for reliable systems capable of accurately and efficiently detecting synthetic media. As deepfakes become more prevalent and easily accessible, developing robust detection methods is crucial to protect privacy, uphold the integrity of digital content, and prevent misuse. This paper aims to contribute to this effort by thoroughly analysing the current state of deepfake detection, supporting the development of more advanced and dependable solutions.

Deepfake technology, a product of advancements in artificial intelligence (AI), specifically deep learning, enables the creation of hyper-realistic synthetic media that can manipulate audio, video, and images to mimic reality convincingly. While this technology offers legitimate applications, such as in entertainment and education, its misuse poses significant societal threats. Deepfakes have been used to spread misinformation, perpetuate fraud, violate individual privacy, and destabilize public trust [1, 2]. The societal implications of deepfake proliferation are profound. For example, deepfakes can undermine democratic processes by creating fabricated political speeches or events [3]. They can also perpetuate personal and institutional damages, such as identity theft and reputation harm [4]. Moreover, the accessibility of deepfake-generating tools exacerbates the problem by enabling individuals with minimal technical expertise to create deceptive content [5]. These issues necessitate urgent attention and robust countermeasures to combat the deepfake menace effectively.

Existing reviews on deepfake technologies primarily focus on foundational concepts and early detection mechanisms. However, the rapid evolution of AI and the growing sophistication of deepfake creation techniques have rendered many of these reviews outdated [6, 7]. This survey aims to fill the gap by providing a comprehensive overview of recent advancements in deepfake detection, prevention, and mitigation strategies. It also emphasizes the importance of addressing the societal and ethical challenges associated with deepfakes [8].

We hypothesize that advancements in machine learning, AI, and cybersecurity offer promising solutions to mitigate the threats posed by deepfakes. By leveraging innovative detection techniques, regulatory frameworks, and collaborative efforts, it is possible to reduce the negative impacts of deepfake technology effectively [9, 10].

This survey is guided by several objectives: to consolidate and evaluate current solutions to the challenges posed by deepfakes, to identify gaps and limitations in existing approaches to deepfake detection and mitigation, and to propose future research directions and strategies for combating deepfake-related issues. The scope of this survey encompasses deepfake detection techniques, including machine learning and digital watermarking methods [2, 9], prevention strategies such as AI-generated content authentication and multi-modal analysis [10], and mitigation efforts, including regulatory frameworks, ethical considerations, and public awareness campaigns [11, 7].

The remainder of this paper is organized as follows: Section 2 reviews deepfake technology, including its evolution, societal implications, and research gaps. Section 3 details the workflow of deepfake detection, highlighting key stages and methodologies. Section 4 outlines detection and mitigation approaches, comparing techniques and evaluation metrics. Section 5 discusses findings, trends, datasets, and mathematical foundations. Section 6 identifies challenges and research gaps, including dataset limitations and real-time detection issues. Section 7 explores recommendations and potential impacts. Section 8 concludes with a summary of findings and the importance of addressing gaps.

2. Literature Review

The proliferation of deepfake technology has prompted extensive research into its origins, advancements, and countermeasures. This section provides a structured review, covering the historical background, key findings, critical analyses, and research gaps in deepfake technology.

2.1. Historical Background

Deepfake technology has transformed the digital landscape, leveraging advancements in artificial intelligence and deep learning. The early foundation of this field was laid with the development of Generative Adversarial Networks (GANs), which facilitated the creation of hyper-realistic visual and audio content [12]. Initially, deepfakes found applications in entertainment and creative industries, such as enhancing visual effects in movies and creating virtual influencers [13]. However, their malicious use for spreading misinformation, violating privacy, and manipulating political narratives has garnered significant attention [14, 15]. The dual-edged nature of this technology highlights both its innovative potential and the ethical dilemmas it poses.

2.2. Key Findings

2.2.1. Categorization of Detection Methods

Research efforts in deepfake detection have yielded several methodologies, each with distinct approaches and objectives:

- **AI-Based Techniques:** Machine learning and deep learning algorithms, particularly Convolutional Neural Networks (CNNs), have achieved notable success in identifying deepfakes by detecting artifacts introduced during the generation process [16, 17]. Advanced models such as recurrent neural networks (RNNs) and transformers have also been explored to analyze temporal inconsistencies in videos [18]. Pre-trained models and transfer learning have further enhanced the efficiency of these techniques.
- **Signal Processing Approaches:** Signal processing-based methods focus on identifying spatial and temporal anomalies in manipulated media. These methods often examine discrepancies in frame transitions, lighting inconsistencies, and unnatural blending between facial regions [19]. Techniques such as spectral analysis and phase correlation are employed to uncover hidden manipulations that are otherwise challenging to detect.
- **Blockchain Solutions:** Blockchain technology is increasingly being adopted for media authentication and traceability. By leveraging immutable ledgers, these solutions can validate the origin and integrity of digital content, thereby providing a robust mechanism to counteract deepfake manipulation [17]. Integration with smart contracts can further automate validation processes, enhancing reliability.
- **Feature Extraction-Based Methods:** Feature extraction-based approaches analyze unique patterns within media to differentiate between authentic and manipulated content. Techniques such as frequency domain analysis, optical flow analysis, and texture-based methods have been employed to identify irregularities that are imperceptible to the human eye [20]. In addition, facial landmark detection and biomechanical consistency checks provide granular insights into potential manipulations.
- **Hybrid Approaches:** Hybrid methods combine multiple techniques, such as integrating AI-based algorithms with signal processing or blockchain frameworks, to enhance detection accuracy. These approaches aim to capitalize on the strengths of each methodology while mitigating their individual limitations [21]. Examples include combining temporal analysis with CNN-based models or integrating blockchain verification with real-time anomaly detection algorithms.

The timeline of deepfake evolution, as shown in Figure 1, provides a detailed overview of the technological advancements that have driven this field. It highlights critical breakthroughs, including the introduction of Generative Adversarial Networks (GANs) in 2014, which revolutionized content generation by enabling high-quality synthetic media. Subsequent developments include advanced autoencoders and transfer learning techniques, which improved model scalability and personalization. The timeline also emphasizes the rise of real-time face reenactment systems, deep neural networks for voice synthesis, and advancements in deepfake detection algorithms. These milestones underline the

rapid growth and sophistication of this technology, posing significant challenges and opportunities in various domains.

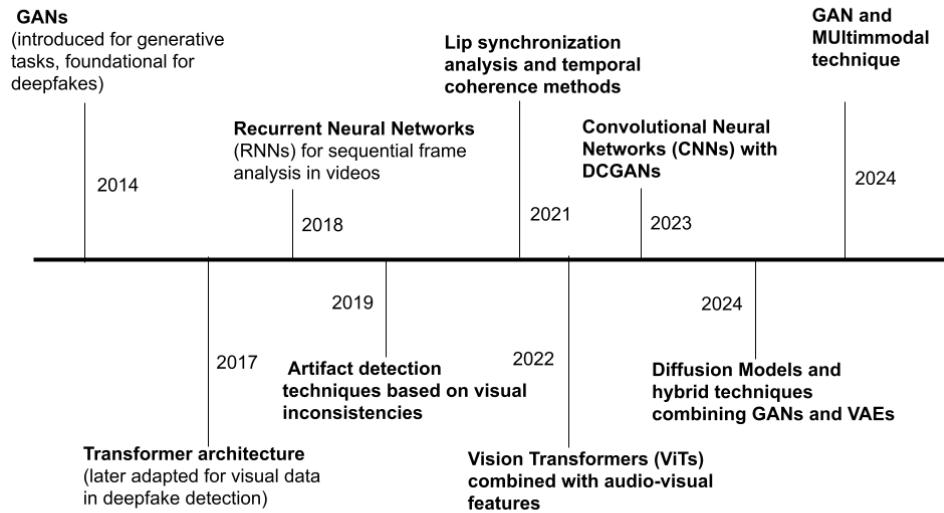


Figure 1: Timeline illustrating the evolution of deepfake technology.

Datasets such as the DeepFake Detection Challenge (DFDC) and FaceForensics++ have underpinned advancements in detection algorithms, providing benchmarks for evaluation [22, 19].

2.3. Critical Analysis

The landscape of deepfake detection is characterized by both significant progress and persistent challenges. AI-driven methods have achieved high accuracy in controlled environments but often struggle with generalization to diverse datasets and unforeseen manipulation techniques [14, 15]. Signal processing approaches, while effective in controlled scenarios, may lack robustness against sophisticated deepfake methods. Blockchain solutions, though promising, face scalability and adoption challenges. Feature extraction techniques are often computationally intensive, limiting their applicability in real-time settings [20, 18].

Recurring issues include the need for standardized evaluation metrics, improved computational efficiency, and ethical considerations. Furthermore, the rapid evolution of deepfake generation methods necessitates continuous adaptation of detection strategies [23, 21]. The absence of datasets that capture real-world variability remains a bottleneck, as most benchmarks are designed for academic purposes [22].

2.4. Identification of Research Gaps

While considerable advancements have been made, several critical gaps remain unaddressed:

- **Real-Time Detection:** The development of lightweight and efficient algorithms capable of real-time processing remains a significant challenge [24]. Advances in edge computing could provide a pathway for achieving this goal.
- **Robustness Across Domains:** Current detection methods require improved generalization to handle diverse datasets and evolving threats [21]. Domain adaptation techniques and unsupervised learning approaches could play a pivotal role.

- **Ethical and Legal Frameworks:** Comprehensive guidelines and regulations addressing the misuse of deepfake technology are urgently needed [25]. Collaboration between technologists, policymakers, and ethicists is essential to establish a robust framework.
- **Advanced Benchmarks:** The lack of standardized and representative datasets hinders the objective evaluation and comparison of detection methods [22]. Future benchmarks should incorporate real-world variations, such as diverse lighting, occlusions, and cultural differences.

Addressing these gaps is imperative for advancing the field of deepfake detection and fostering trust in digital ecosystems. Future research must prioritize the development of scalable, robust, and ethically aligned solutions to counteract the growing threats posed by deepfake technology. Collaboration across disciplines and the integration of emerging technologies will be key to overcoming these challenges.

3. Workflow: Deepfake Detection

The process of deepfake detection involves several critical steps, as illustrated in Figure 2. Each step plays a vital role in accurately distinguishing between original and fake content. Below is a detailed explanation of the workflow, along with examples of methodologies and techniques commonly employed at each stage:

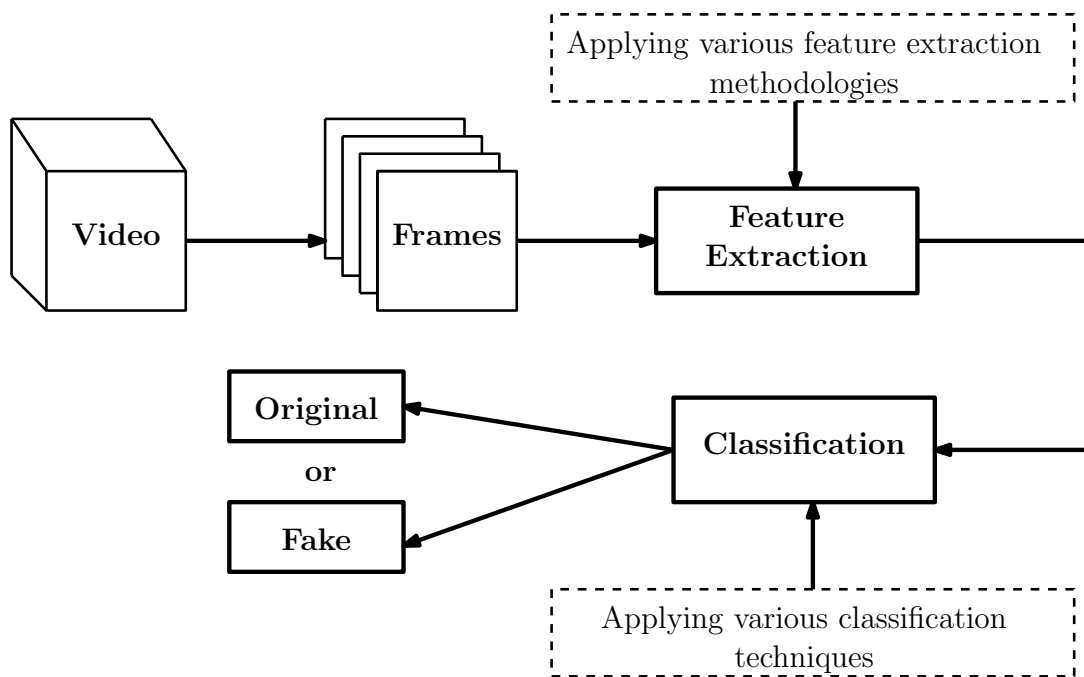


Figure 2: Workflow illustrating the steps in deepfake detection.

1. Input - Video Frames Extraction

The first step involves splitting the input video into individual frames. These frames serve as the foundational data for further analysis. High-resolution frames are preferred to ensure the features used in detection are well-represented.

Example Methodology:

- **Frame Sampling:** Extract frames at fixed intervals (e.g., every nth frame) to reduce computational load while maintaining key details.

2. Feature Extraction

Feature extraction involves identifying and isolating the most critical aspects of the video frames that can reveal inconsistencies or unnatural patterns. These features form the basis for differentiating between real and fake media.

Example Feature Extraction Methods:

- **Pixel-Level Artifacts Detection:** Focus on artifacts such as inconsistent lighting, shadows, or pixel distortions often introduced during deepfake generation.
- **Temporal Inconsistencies:** Analyze frame-to-frame transitions for unnatural movement or discontinuities.
- **Frequency Domain Analysis:** Techniques like Discrete Fourier Transform (DFT) or Wavelet Transform to detect anomalies in high-frequency bands.
- **Biometric Feature Analysis:** Focus on facial landmarks, eye movement, and lip-sync patterns to identify irregularities.

3. Classification

Once features are extracted, they are fed into a classification model to predict whether the content is original or fake. This step leverages machine learning and deep learning algorithms to make the final determination.

Example Classification Techniques:

- **Traditional Machine Learning Models:**
 - **SVM (Support Vector Machines):** Effective for small datasets and well-defined features.
 - **Random Forest:** Ensemble-based approach for feature importance and classification.
- **Deep Learning Models:**
 - **Convolutional Neural Networks (CNNs):** Suitable for spatial features like pixel-level inconsistencies or facial biometrics.
 - **Recurrent Neural Networks (RNNs):** Ideal for temporal features such as frame continuity and motion consistency.
 - **EfficientNet, MobileNetV2, and VGG16:** Pretrained architectures fine-tuned for deepfake detection tasks.
- **Hybrid Models:** Combining CNNs for spatial features with RNNs for temporal consistency checks.

4. Output - Classification Result

The final step produces a classification result that labels the input as either "**Original**" or "**Fake**". The accuracy and reliability of this output depend on the effectiveness of the previous steps and the quality of training data used to build the detection model.

Evaluation Metrics:

- **Accuracy, Precision, Recall:** Measure overall model performance.
- **F1 Score:** Balance between precision and recall.
- **AUC-ROC Curve:** Evaluate model sensitivity to different thresholds.

4. Methodologies and Approaches

This section outlines the methodologies employed in surveying the research landscape on deepfake detection and mitigation. It describes the survey methodology, provides detailed insights into various approaches analyzed, and presents a comparative analysis of these methodologies.

4.1. Survey Methodology

The reviewed papers were selected using a systematic approach to ensure comprehensive coverage of the field. A database search was conducted across platforms such as IEEE Xplore, Springer, and ACM Digital Library using keywords like "deepfake detection," "GAN-based manipulation," and "blockchain authentication." The inclusion criteria prioritized articles published in peer-reviewed journals and conferences between 2019 and 2025. Studies that lacked empirical results or focused solely on deepfake generation without discussing detection were excluded. A total of 50 papers met these criteria and were included in this review.

The evaluation framework for categorizing existing solutions focused on three key dimensions:

- **Technique:** Classification into AI-based, signal processing-based, blockchain-assisted, handcrafted feature extraction, and hybrid approaches.
- **Performance Metrics:** Accuracy, scalability, and computational efficiency.
- **Applicability:** Suitability for real-time detection and generalization across datasets.

4.2. Approaches Analyzed

4.2.1. Machine Learning/AI-Based Techniques

Machine learning and AI-based techniques are among the most widely explored methods for deepfake detection. Convolutional Neural Networks (CNNs) effectively detect spatial inconsistencies, such as unnatural textures and blending artifacts, in manipulated media [16]. Recurrent Neural Networks (RNNs) and transformers analyze temporal patterns, making them well-suited for video analysis [18]. Generative Adversarial Networks (GANs), while primarily used for creating deepfakes, are also utilized for adversarial training to identify and counteract synthetic content [20]. Furthermore, pre-trained models and transfer learning approaches have improved detection performance by reducing training requirements and leveraging pre-existing knowledge bases.

4.2.2. Digital Forensics Techniques

Digital forensics relies on analyzing inconsistencies and artifacts in video and audio signals. Techniques such as phase correlation, frequency domain analysis, and optical flow detection identify discrepancies that are challenging for deepfake generation algorithms to mimic [19]. For instance, variations in lighting, unnatural reflections, and irregularities in motion provide telltale signs of manipulation. These methods are particularly valuable in scenarios where content integrity is under scrutiny.

4.2.3. Blockchain for Authentication

Blockchain technology offers a robust framework for verifying the authenticity and provenance of digital content. Immutable ledgers record the history of media, ensuring traceability and preventing tampering [17]. Smart contracts enable automated verification processes, enhancing the scalability of blockchain-assisted solutions. This approach is particularly effective in applications requiring real-time validation, such as social media and news dissemination platforms.

4.2.4. Handcrafted Feature Extraction Techniques

Handcrafted feature extraction focuses on identifying specific features that distinguish manipulated from authentic media. These methods analyze elements such as facial landmarks, eye blinking patterns, and lip synchronization [20]. Techniques like Local Binary Patterns (LBP) and Histogram of Oriented Gradients (HOG) are used to detect texture inconsistencies and unnatural movements. Although computationally less intensive than AI-based approaches, handcrafted techniques often struggle with the subtle sophistication of modern deepfakes.

4.2.5. Hybrid Approaches

Hybrid approaches integrate multiple methodologies to enhance robustness and accuracy. For example, combining CNNs with optical flow analysis leverages both spatial and temporal insights. Similarly, blockchain verification can be paired with AI-driven anomaly detection for comprehensive validation [21]. These approaches aim to balance the strengths of individual techniques while mitigating their limitations, making them suitable for complex and diverse use cases.

4.3. Comparative Analysis

A comparative analysis of the methodologies is presented in Table 1, highlighting their efficiency, accuracy, scalability, and suitability for real-time detection.

Table 1
Comparison of Deepfake Detection Approaches

Approach	Accuracy	Scalability	Real-Time Suitability
AI-Based Techniques	High	Moderate	Limited due to computational intensity
Digital Forensics	Moderate	High	Moderate
Blockchain Solutions	High	Low	High
Handcrafted Feature Extraction	Moderate	High	High
Hybrid Approaches	Very High	Moderate	High

AI-based techniques excel in accuracy but are computationally demanding, making scalability and real-time application challenging. Digital forensics methods offer high scalability but may struggle with sophisticated manipulations. Blockchain solutions provide high reliability and real-time suitability but face scalability issues due to resource requirements. Handcrafted feature extraction methods are efficient and scalable but less effective against subtle manipulations. Hybrid approaches represent a balanced solution, combining accuracy, scalability, and real-time suitability.

In conclusion, while each methodology has its strengths and weaknesses, hybrid approaches demonstrate the most promise for addressing the diverse challenges posed by deepfake technology.

5. Findings and Trends

5.1. Key Insights

Recent advancements in deepfake detection have introduced innovative techniques that significantly improve accuracy and robustness against increasingly sophisticated deepfake content. Maheshwari et al. (2024) explored plasmonic nanomaterials with surface plasmon resonance (SPR) for image detection, achieving over 95% accuracy even in complex scenarios [26]. A hybrid deep learning model combining MesoNet4 and ResNet101 was proposed by Javed et al. (2024), attaining detection accuracies of 98.73%, 96.89%, and 97.90% on FaceForensics++, CelebV1, and CelebV2 datasets, respectively [27]. Advanced biosensors integrating plasmonic resonance with convolutional neural networks reached 98.7% accuracy and demonstrated rapid response times (0.8 seconds per frame) [28].

Blockchain-based federated learning approaches, such as Heidari et al.'s (2024) method, enhanced accuracy by 6.6% compared to benchmarks while maintaining data confidentiality [29]. Temporal feature prediction schemes focusing on audio-visual modalities demonstrated superior accuracy (84.33%) on the FakeAVCeleb dataset [30]. Vision Transformers (ViTs) showed great promise in multiclass detection tasks, achieving an F1-score of 99.90%, outperforming traditional CNNs [31]. Kingra et al.'s (2024) SFormer architecture, based on spatio-temporal transformers, achieved up to 100% accuracy on datasets such as FF++ and Deeper-Forensics [32]. Almetekawy et al. (2024) demonstrated that incorporating

spatiotemporal textures improved reproducibility and accuracy by up to 91.96% [33]. Guarnera et al. (2024) introduced a hierarchical multi-level approach for deepfake detection, achieving 97% accuracy across multiple GAN and diffusion model tasks [34]. Gao et al. (2024) used temporal audio-video feature prediction to reach an 84.33% accuracy [30]. Lastly, Arshed et al. (2024) explored Vision Transformers (ViTs) achieving F1-scores close to 99.90% [31].

5.2. Statistical Analysis

Table 2 compares the performance metrics, including accuracy, computational cost, and dataset benchmarks, for different methods. These approaches reflect varying trade-offs in sensitivity, speed, and dataset applicability.

Table 2
Performance Comparison of Deepfake Detection Techniques

Technique	Accuracy	Dataset(s)	Strengths/Weaknesses
Plasmonic Nanomaterials (SPR) [26]	95%	Custom Dataset	High sensitivity; robust to lighting conditions but computationally intensive.
Hybrid Model (MesoNet4 + ResNet101) [27]	98.73% (FaceForensics++), 96.89% (CelebV1)	FaceForensics++, CelebV1, CelebV2	Real-time capability; limited multimodal application.
Advanced Plasmonic Biosensor [28]	98.7%	Custom dataset	Fast response time; integration challenges in real-world scenarios.
Blockchain-Based Federated Learning [29]	6.6% increase over benchmarks	Diverse	Data privacy maintained; high computational complexity.
Temporal Feature Prediction [30]	84.33%	FakeAVCeleb	Novel audio-visual fusion; lower accuracy compared to transformer-based models.
Vision Transformers (ViTs) [31]	99.90%	Multiclass-prepared dataset	High accuracy; robust to compression and resizing.
SFormer [32]	Up to 100%	FF++, Deeper-Forensics	Superior generalization; computationally expensive.
Spatiotemporal Textures [33]	91.96%	Celeb-DF, FF++	Enhanced stability; moderate accuracy in cross-dataset scenarios.
Hierarchical Multi-level GAN Analysis [34]	97%	GAN and Diffusion Model Dataset	Robust to attacks like compression; lacks real-time capabilities.
Patch-Wise Deep Learning [31]	99.90%	GAN, Stable Diffusion Datasets	Impressive F1 rates; computational overhead.

5.3. Popular Datasets for Deepfake Validation

Datasets play a crucial role in validating and improving deepfake detection solutions. Table 3 highlights some of the most popular datasets used in this domain, emphasizing their size, types of content, and primary applications.

5.4. Emerging Trends

Several trends in deepfake detection have emerged:

- **Multimodal Solutions:** Techniques like temporal feature prediction and hybrid models increasingly integrate multiple modalities (e.g., audio and video) to enhance detection accuracy

Table 3
Popular Datasets Used for Deepfake Validation

Dataset Name	Description	Size	Types of Content	Source
FaceForensics++	Large-scale dataset for face forgery detection.	1,000 videos	Deepfake, Neural Rendered, Face Swapping	University of Erlangen-Nuremberg
DeepFakeDetection	Focused on detecting deepfake videos.	3,000 videos	Real, Deepfake	University of California, Berkeley
Celeb-DF	High-resolution deepfake videos featuring celebrities.	5,639 videos	Celebrities, TV Shows	Zhejiang University
DFDC (Deepfake Detection Challenge)	Comprehensive dataset for deepfake challenges.	100,000 videos	Real, Deepfake	Facebook AI
The Realism of Deepfakes	Evaluates realism in deepfake generation.	Fully Annotated	Deepfake, GANs	Stanford University

[30].

- **Transformer Architectures:** Vision Transformers (ViTs) and spatio-temporal transformer models like SFormer demonstrate exceptional performance, particularly in generalizing across datasets [31, 32].
- **Adversarial Learning:** GAN-based methods for deepfake generation have inspired adversarial learning approaches to detect increasingly realistic fakes.
- **Real-Time and Scalable Solutions:** Biosensors and hybrid architectures focus on reducing latency, with potential for real-time applications [28, 27].
- **Privacy-Preserving Techniques:** Blockchain-based federated learning represents a shift towards safeguarding data privacy while achieving robust detection [29].

These trends indicate a paradigm shift towards integrating diverse modalities, leveraging advanced architectures, and prioritizing real-time and privacy-preserving solutions for scalable and effective deepfake detection.

5.5. Mathematical Foundations for Detection

In deepfake detection, various mathematical models and techniques are employed to enhance accuracy and robustness. The key mathematical foundations for these detection models include Generative Adversarial Networks (GANs), Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Attention Mechanisms, Adversarial Training Loss, and Ensemble Prediction. As shown in Table 4, GANs leverage an adversarial training approach, where a generator and discriminator interact to distinguish real from fake data. CNNs, on the other hand, apply convolution operations to extract spatial features from images, crucial for analyzing image patterns in deepfakes. RNNs are employed for sequential data, such as video frames, to capture temporal dependencies. The attention mechanism, often used in Vision Transformers (ViTs), helps models focus on significant features, enhancing the detection process. Additionally, adversarial training loss is designed to improve model robustness by exposing it to adversarial examples. Finally, ensemble prediction aggregates results from multiple models to boost the overall detection accuracy.

6. Challenges and Gaps

6.1. Current Challenges

Despite the advancements in deepfake detection technologies, several technical challenges persist, limiting the effectiveness of current solutions:

Table 4
Mathematical Formulas for Deepfake Detection Models

Mathematical Concept	Formula and Explanation
1. Generative Adversarial Networks (GANs)	$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$ <p> $D(x)$: Discriminator's probability of x being real. $G(z)$: Data generated by G from noise z. $p_{\text{data}}(x)$: Distribution of real data. $p_z(z)$: Distribution of noise. </p>
2. Convolutional Neural Networks (CNNs)	$y[i, j] = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} x[i+m, j+n] \cdot k[m, n]$ <p> $x[i, j]$: Input image pixel at position (i, j). $k[m, n]$: Filter kernel of size $M \times N$. $y[i, j]$: Convolved output. </p>
3. Recurrent Neural Networks (RNNs)	$h_t = \sigma(W_h h_{t-1} + W_x x_t + b)$ <p> h_t: Hidden state at time t. h_{t-1}: Hidden state from the previous time step. x_t: Input at time t. W_h, W_x: Weight matrices. b: Bias vector. σ: Activation function. </p>
4. Attention Mechanism	$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$ <p> Q, K, V: Query, key, and value matrices. d_k: Dimension of the key vector. </p>
5. Adversarial Training Loss	$\mathcal{L}_{\text{adv}} = \mathbb{E}_{(x, y) \sim \mathcal{D}} [\max_{\delta \in S} \ell(f(x + \delta), y)]$ <p> δ: Perturbation within constraint S. $\ell(f(x), y)$: Loss function comparing prediction $f(x)$ with label y. </p>
6. Ensemble Prediction	$P_{\text{ensemble}} = \frac{1}{N} \sum_{i=1}^N P_i$ <p> P_i: Prediction probability from the i-th model. N: Number of models. </p>

- **Detection Accuracy for Low-Quality Videos:** Many deepfake detection models struggle with low-resolution or highly compressed videos, which are often encountered on social media platforms. This degradation in quality obscures telltale artifacts, reducing detection performance.
- **Computational Overhead:** Deep learning-based detection methods, while highly accurate, often require significant computational resources. Balancing the need for high detection accuracy with computational efficiency remains a key challenge, particularly for real-time applications.
- **Generalization Across Techniques:** As new and more sophisticated deepfake generation techniques emerge, detection models often fail to generalize, requiring constant retraining on updated datasets.
- **Real-Time Detection:** Many existing approaches lack the speed needed for real-time detection, especially in live-streaming or high-throughput environments, where immediate detection is crucial.
- **Robustness to Adversarial Attacks:** Deepfake detection models are vulnerable to adversarial attacks that subtly alter fake content to evade detection mechanisms.

6.2. Research Gaps

In addition to technical challenges, there are several gaps in current research that must be addressed to advance deepfake detection methodologies:

- **Standardized Datasets:** While several datasets exist, there is a lack of universally accepted

benchmarks that cover diverse content types, resolutions, and manipulation techniques. Creating standardized, diverse datasets would enhance model comparability and reliability.

- **Legal and Ethical Frameworks:** Deepfake detection research often overlooks the legal and ethical implications of using synthetic media. Establishing guidelines for the responsible use of detection technologies and addressing privacy concerns is critical.
- **Robustness Against Evolving Deepfake Techniques:** As generative models continue to evolve, there is a need for detection methods that can adapt to new manipulation techniques without requiring frequent retraining.
- **Cross-Platform Scalability:** Detection methods often perform well on specific datasets but fail when deployed across different platforms or real-world scenarios. Research into scalable and robust cross-platform solutions is necessary.
- **Human-AI Collaboration:** Current systems primarily focus on automated detection, with little emphasis on integrating human expertise to improve accuracy and interpretability of results.
- **Ethical Use of Detection Tools:** There is a need to address potential misuse of detection tools themselves, such as leveraging them to create more advanced deepfakes by understanding their weaknesses.

Addressing these challenges and research gaps will require a concerted effort from academia, industry, and policymakers to ensure that deepfake detection technologies remain effective, equitable, and ethical in the face of evolving threats.

7. Future Directions

7.1. Recommendations

To advance the field of deepfake detection and mitigate the risks associated with synthetic media, the following actionable steps are recommended:

- **Development of Lightweight, Real-Time Models:** Future research should focus on creating computationally efficient deepfake detection models capable of real-time processing. This involves exploring novel architectures, such as transformer-based models optimized for speed and scalability.
- **Building More Diverse and Representative Datasets:** Establishing datasets that include a wide variety of manipulation techniques, demographics, and content types will improve the robustness and generalizability of detection models. Collaboration among research institutions and industry can facilitate the creation of comprehensive benchmarks.
- **Creating Legal and Ethical Frameworks:** Policymakers and researchers should work together to establish guidelines for the responsible use of generative technologies. This includes defining acceptable practices, ensuring transparency, and addressing privacy concerns in dataset usage.
- **Enhancing Robustness Against Adversarial Attacks:** Research should prioritize techniques to make detection models resilient to adversarial examples, such as adversarial training, ensemble methods, and anomaly detection frameworks.
- **Integration of Multimodal Approaches:** Combining audio, video, and textual data can lead to more comprehensive detection systems. Future work should focus on integrating these modalities effectively to improve detection accuracy.
- **Fostering Human-AI Collaboration:** Developing tools that allow human experts to interact with detection systems can enhance the interpretability and reliability of results, particularly in high-stakes scenarios.

7.2. Potential Impact

The proposed advancements in deepfake detection can have far-reaching implications across various domains:

- **Policy-Making:** Improved detection methods and standardized datasets can inform regulatory frameworks, helping governments and organizations address the ethical and legal challenges posed by deepfake technology.
- **Societal Trust:** By effectively mitigating the spread of synthetic media, advanced detection technologies can restore public trust in digital content, reducing the impact of misinformation and manipulation.
- **Adoption of AI Technologies:** The development of robust and ethical deepfake detection systems will encourage the responsible adoption of AI technologies in industries such as media, entertainment, and cybersecurity.
- **Enhanced Security Measures:** Real-time detection capabilities can be integrated into digital platforms, safeguarding users against malicious deepfake content and protecting sensitive information.

By addressing these recommendations and leveraging the potential impact, the research community can ensure that deepfake detection technologies remain a step ahead of evolving generative methods, fostering a safer and more trustworthy digital environment.

8. Conclusion

This survey has explored the current state of deepfake detection technologies, highlighting the rapid advancements in methods designed to counteract the growing sophistication of generative models. Key insights include the effectiveness of hybrid approaches, such as combining multimodal analysis with AI-based techniques, and the potential of transformer-based architectures to improve accuracy and scalability. Despite these advancements, challenges persist in detecting low-quality or adversarially manipulated deepfakes, underscoring the need for robust and adaptable solutions.

This work consolidates knowledge from diverse fields, presenting a comprehensive review of the strengths and limitations of existing deepfake detection methods. By identifying research gaps—such as the need for standardized datasets and ethical frameworks—this survey provides a roadmap for future studies. It also emphasizes the importance of integrating human expertise with automated systems to enhance the interpretability and reliability of detection outcomes.

As deepfake technology continues to evolve, the importance of proactive research and collaboration cannot be overstated. The development of lightweight, real-time detection models and the establishment of legal and ethical standards are crucial steps toward combating the misuse of synthetic media. By fostering cross-disciplinary partnerships and prioritizing innovation, the research community can address emerging threats and ensure the responsible use of AI technologies, safeguarding societal trust and digital integrity.

Declaration on Generative AI

The author(s) have not employed any Generative AI tools.

References

- [1] S. H. Al-Khazraji, H. H. Saleh, A. I. KHALID, I. A. MISHKHAL, Impact of deepfake technology on social media: Detection, misinformation and societal implications, *The Eurasia Proceedings of Science Technology Engineering and Mathematics* 23 (2023) 429–441.
- [2] M. Sharma, M. Kaur, A review of deepfake technology: an emerging ai threat, *Soft Computing for Security Applications: Proceedings of ICSCS 2021* (2022) 605–619.
- [3] C. Whyte, Deepfake news: Ai-enabled disinformation as a multi-level public policy challenge, *Journal of cyber policy* 5 (2020) 199–217.

- [4] P. Singh, D. B. Dhiman, Exploding ai-generated deepfakes and misinformation: A threat to global concern in the 21st century, Available at SSRN 4651093 (2023).
- [5] W. Matli, Extending the theory of information poverty to deepfake technology, *International Journal of Information Management Data Insights* 4 (2024) 100286.
- [6] A. O. Kwok, S. G. Koh, Deepfake: a social construction of technology perspective, *Current Issues in Tourism* 24 (2021) 1798–1802.
- [7] D. Chapagain, N. Kshetri, B. Aryal, Deepfake disasters: A comprehensive review of technology, ethical concerns, countermeasures, and societal implications, in: *2024 International Conference on Emerging Trends in Networks and Computer Communications (ETNCC)*, IEEE, 2024, pp. 1–9.
- [8] D. Sarkar, S. De Sarkar, Combatting deep-fakes in india—an analysis of the evolving legal paradigm and its challenges (2024).
- [9] M. Mustak, J. Salminen, M. Mäntymäki, A. Rahman, Y. K. Dwivedi, Deepfakes: Deceptions, mitigations, and opportunities, *Journal of Business Research* 154 (2023) 113368.
- [10] M. R. Shoaib, Z. Wang, M. T. Ahvanooy, J. Zhao, Deepfakes, misinformation, and disinformation in the era of frontier ai, generative ai, and large ai models, in: *2023 International Conference on Computer and Applications (ICCA)*, IEEE, 2023, pp. 1–7.
- [11] M. Pawelec, Decent deepfakes? professional deepfake developers’ ethical considerations and their governance potential, *AI and Ethics* (2024) 1–26.
- [12] R. Chataut, A. Upadhyay, Introduction to deepfake technology and its early foundations, in: *Deepfakes and Their Impact on Business*, IGI Global Scientific Publishing, 2025, pp. 1–18.
- [13] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, J. Ortega-Garcia, Deepfakes and beyond: A survey of face manipulation and fake detection, *Information Fusion* 64 (2020) 131–148.
- [14] P. Yu, Z. Xia, J. Fei, Y. Lu, A survey on deepfake video detection, *Iet Biometrics* 10 (2021) 607–624.
- [15] A. Malik, M. Kuribayashi, S. M. Abdullahi, A. N. Khan, Deepfake detection for human face images and videos: A survey, *Ieee Access* 10 (2022) 18757–18775.
- [16] A. Heidari, N. Jafari Navimipour, H. Dag, M. Unal, Deepfake detection using deep learning methods: A systematic and comprehensive review, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 14 (2024) e1520.
- [17] M. S. Rana, M. N. Nobli, B. Murali, A. H. Sung, Deepfake detection: A systematic literature review, *IEEE access* 10 (2022) 25494–25513.
- [18] B. Kaddar, S. A. Fezza, Z. Akhtar, W. Hamidouche, A. Hadid, J. Serra-Sagrìstà, Deepfake detection using spatiotemporal transformer, *ACM Transactions on Multimedia Computing, Communications and Applications* (2024).
- [19] Y. Nirkin, L. Wolf, Y. Keller, T. Hassner, Deepfake detection based on discrepancies between faces and their context, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44 (2021) 6111–6121.
- [20] T. Zhang, Deepfake generation and detection, a survey, *Multimedia Tools and Applications* 81 (2022) 6259–6276.
- [21] Y. Patel, S. Tanwar, R. Gupta, P. Bhattacharya, I. E. Davidson, R. Nyameko, S. Aluvala, V. Vimal, Deepfake generation and detection: Case study and challenges, *IEEE Access* (2023).
- [22] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, C. C. Ferrer, The deepfake detection challenge (dfdc) dataset, *arXiv preprint arXiv:2006.07397* (2020).
- [23] J. W. Seow, M. K. Lim, R. C. Phan, J. K. Liu, A comprehensive overview of deepfake: Generation, detection, datasets, and opportunities, *Neurocomputing* 513 (2022) 351–371.
- [24] M. Weerawardana, T. Fernando, Deepfakes detection methods: A literature survey, in: *2021 10th International Conference on Information and Automation for Sustainability (ICIAfS)*, IEEE, 2021, pp. 76–81.
- [25] A. Chadha, V. Kumar, S. Kashyap, M. Gupta, Deepfake: an overview, in: *Proceedings of second international conference on computing, communications, and cyber-security: IC4S 2020*, Springer, 2021, pp. 557–566.
- [26] R. U. Maheshwari, B. Paulchamy, B. K. Pandey, D. Pandey, Enhancing sensing and imaging capabilities through surface plasmon resonance for deepfake image detection, *Plasmonics* (2024)

1–20.

- [27] M. Javed, Z. Zhang, F. H. Dahri, A. A. Laghari, Real-time deepfake video detection using eye movement analysis with a hybrid deep learning approach, *Electronics* 13 (2024) 2947.
- [28] R. U. Maheshwari, S. Kumarganesh, S. KVM, A. Gopalakrishnan, K. Selvi, B. Paulchamy, P. Rishabavarthani, K. M. Sagayam, B. K. Pandey, D. Pandey, Advanced plasmonic resonance-enhanced biosensor for comprehensive real-time detection and analysis of deepfake content, *Plasmonics* (2024) 1–18.
- [29] A. Heidari, N. J. Navimipour, H. Dag, S. Talebi, M. Unal, A novel blockchain-based deepfake detection method using federated and deep learning models, *Cognitive Computation* (2024) 1–19.
- [30] Y. Gao, X. Wang, Y. Zhang, P. Zeng, Y. Ma, Temporal feature prediction in audio–visual deepfake detection, *Electronics* 13 (2024) 3433.
- [31] M. A. Arshed, S. Mumtaz, M. Ibrahim, C. Dewi, M. Tanveer, S. Ahmed, Multiclass ai-generated deepfake face detection using patch-wise deep learning model, *Computers* 13 (2024) 31.
- [32] S. Kingra, N. Aggarwal, N. Kaur, Sformer: An end-to-end spatio-temporal transformer architecture for deepfake detection, *Forensic Science International: Digital Investigation* 51 (2024) 301817.
- [33] A. Almestekawy, H. H. Zayed, A. Taha, Deepfake detection: Enhancing performance with spatiotemporal texture and deep learning feature fusion, *Egyptian Informatics Journal* 27 (2024) 100535.
- [34] L. Guarnera, O. Giudice, S. Battiato, Level up the deepfake detection: a method to effectively discriminate images generated by gan architectures and diffusion models, in: *Intelligent Systems Conference*, Springer, 2024, pp. 615–625.